

Version with Markings to Show Changes Made

1 1. A method for managing caches in a system with multiple caches that may contain
2 different copies of a data item, comprising the steps of:
3 modifying the data item in a first node of said multiple caches to create a modified
4 data item;
5 sending the modified data item from said first node to a second node of said multiple
6 caches without durably storing the modified data item from said first node to
7 persistent storage;
8 after said modified data item has been sent from said first node to said second node,
9 said first node sending a request to a master of said data item for writing said
10 data item to persistent storage; and
11 in response to said request, said master coordinating with said multiple caches to
12 cause said data item to be written to persistent storage.

1 2. The method of Claim 1 wherein:
2 the method includes the step of maintaining, within an ordered series of bins, entries
3 for past-image versions of data items;
4 each bin in said ordered series corresponds to a time range;
5 a particular bin corresponds to the time range that covers the time at which the data
6 item is modified in said first node; and
7 the step of sending a request is performed by sending a request for writing a particular
8 bin of said ordered series of bins to persistent storage.

1 3. The method of Claim 2 wherein the step of said master coordinating with said
2 multiple caches to cause said data item to be written to persistent storage includes
3 said master causing said multiple caches to write data items to persistent storage to

4 cover all past image versions of data items that were modified during the time range
5 of said particular bin.

1 4. The method of Claim 3 further comprising the step of emptying said particular bin
2 after said multiple caches write data items to persistent storage to cover all past image
3 versions of data items that were modified during the time range of said particular bin.

1 5. The method of Claim 4 wherein the step of emptying said particular bin includes the
2 steps of:
3 discarding entries within said particular bin that are associated with past images that
4 have last-dirtied times within the time range of said particular bin; and
5 moving to one or more other bins the entries within said particular bin that are
6 associated with past images that have last-dirtied times later than the time
7 range of said particular bin.

1 6. The method of Claim 1 wherein the step of sending a request to a master is performed
2 by sending the request to a global lock manager.

1 7. The method of Claim 1 wherein the step of sending a request to a master is performed
2 by sending the request to a lock manager that is one of a plurality of lock managers
3 within a distributed lock management system.

1 8. The method of Claim 1 further comprising the step of sending from the master, to
2 interested nodes, write-notification messages indicating that said data item has been
3 written to persistent storage, in response to said data item being written to persistent
4 storage.

- 1 9. The method of Claim 8 wherein the step of sending write-notification messages
2 includes the master sending to at least one interested node a single message that
3 notifies said at least one interested node that a plurality of data items have been
4 written to persistent storage.
- 1 10. The method of Claim 1 wherein the step of said first node sending a request to a
2 master of said data item for writing said data item to persistent storage includes the
3 first node sending to said master a single message that requests writing a plurality of
4 data items to persistent storage, wherein said plurality of data items includes said data
5 item.
- 1 11. The method of Claim 10 wherein the step of sending a single message includes
2 sending a message that identifies a bin to request that all data items that belong to the
3 bin be written to persistent storage.
- 1 12. The method of Claim 11 wherein the bin is associated with a range of time and
2 includes data items that were first dirtied by the first node during said range of time
3 and that were subsequently transferred to other nodes without first being written to
4 persistent storage.
- 1 13. The method of Claim 8 wherein the step of sending from the master to interested
2 nodes write-notification messages includes the steps of:
3 immediately sending write-notification messages to a first set of interested nodes,
4 where said first set of interested nodes includes the interested nodes that have
5 requested said data item to be written to persistent storage; and

6 delaying the sending of write-notification messages to a second set of nodes, where
7 said second set of nodes includes interested nodes that do not belong to said
8 first set of interested nodes.

1 14. The method of Claim 8 wherein the step of sending from the master to interested
2 nodes includes delaying the sending of write-notification messages to at least one
3 interested node.

1 15. The method of Claim 14 wherein a write-notification message is sent to the at least
2 one interested node in response to a lock request made by said at least one interested
3 node.

1 16. The method of Claim 14 wherein a write-notification message is sent to the at least
2 one interested node in response to the at least one interested node requesting that said
3 data item be written to persistent storage.

1 17. The method of Claim 14 wherein a write-notification message is sent to the at least
2 one interested node within a ping request that the master sends to the at least one
3 interested node for the at least one interested node to transfer another data item to
4 another node.

1 18. The method of Claim 1 wherein the step of coordinating includes the steps of:
2 determining whether a version of said data item, that is at least as recent as said
3 modified version, has already been written to persistent storage; and
4 if a version of said data item that is at least as recent as said modified version has
5 already been written to persistent storage, then sending a write-notification
6 message from said master to notify said first node that a version of said data

7 item that is at least as recent as said modified version has already been written
8 to persistent storage.

1 19. The method of Claim 18 wherein the step of coordinating includes, if a version of
2 said data item that is at least as recent as said modified version has not already been
3 written to persistent storage, then sending a write-perform message from said master
4 to grant permission for said modified version to be written to persistent storage.

1 20. The method of Claim 1 wherein the step of coordinating includes the
2 steps of:
3 selecting a particular node of said multiple caches that has a particular version of said
4 data item, wherein said particular version is at least as recent as the modified
5 data item in said first node; and
6 causing said particular version of said data item to be written from said particular
7 node to persistent storage.

1 21. The method of Claim 20 wherein the step of selecting a particular node includes
2 selecting the node, of said multiple caches, that has a most recently modified version
3 of said data item.

1 22. The method of Claim 20 further comprising the step of the master informing the first
2 node that said data item has been written to persistent storage in response to the
3 master receiving confirmation that said particular version of said data item has been
4 written to persistent storage.

1 23. The method of Claim 20 further comprising the step of the master informing a set of
2 caches that said data item has been written to persistent storage in response to the
3 master receiving confirmation that said particular version of said data item has been

written to persistent storage, wherein said set of caches includes caches, other than said particular node, that contain modified versions of said data item that are not more recent than said particular version.

24. A method for managing caches in a system with multiple caches that may contain different copies of a data item, comprising the steps of:
modifying the data item in a first cache to create a modified data item;
in response to writing the modified data item to persistent storage, performing the steps of:
a node associated with the first cache determining whether any other cache in said multiple caches had created a dirty version of said data item; and
if any other cache in said multiple caches had created a dirty version of said data item, then the node associated with the first cache informing a master of said data item that said modified data item has been written to persistent storage; and
if no other cache in said multiple caches had created a dirty version of said data item, then the step of writing the modified data item is performed without informing said master that said modified data item has been written to persistent storage.

25. The method of Claim 24 wherein the step of determining whether any other cache in said multiple caches had created a dirty version of said data item includes inspecting a global dirty flag associated with said data item.

26. The method of Claim 25 wherein:
prior to modifying the data item in said first cache, the data item had been modified in a second cache;

the data item is not persistently stored between being modified in said second cache
and being modified in said first cache; and
a node associated with the second cache causes said global dirty flag to be set to
indicate that said data item is globally dirty.

27. The method of Claim 26 wherein the node associated with the second cache causes
said global dirty flag to be set in response to transferring said data item from said
second cache to another cache of said multiple caches.

28. A method for managing caches in a system with multiple caches that may contain
different copies of a data item, comprising the steps of:
modifying the data item in a first cache to create a modified data item;
when a node associated with the first node desires to write said modified data item to
persistent storage, performing the steps of:
if then node associated with the first cache does not currently have ownership
rights to said data item, then the node associated with the first cache
sending a request to a master of said data item for said data item to be
written to persistent storage; and
if said node associated with said first cache currently has ownership rights to
said data item, then the first node writing said data item to persistent
storage without sending a request to said master for said data item to
be written to persistent storage.

~~28~~29. The method of Claim ~~27~~28 further comprising the step of designating holders of
exclusive locks in data items to be owners of said data items.

~~29~~30. A method for managing a data item, the method comprising the steps of:

2 when a node that has an exclusive lock on a data item desires to write the data item to
3 persistent storage, performing the steps of
4 determining whether a mode associated with the data item is local or global;
5 if the mode associated with the data item is local, then the node writing the
6 data item to persistent storage without communicating with a master of
7 said data item; and
8 if the mode associated with the data item is global, then the node sending a
9 message to the master of the data item to request writing of said data
10 item to persistent storage.

1 ~~30~~31. The method of Claim ~~29~~30 wherein:

2 the mode associated with the data item is global; and

3 the method further includes:

4 the node receiving permission from the master to write the data item to
5 persistent storage; and

6 after writing the data item to persistent storage, changing the mode from
7 global to local.

1 ~~31~~32. The method of Claim ~~29~~30 wherein:

2 the mode associated with the data item is local; and

3 before the node has completed writing of the data item to persistent storage, the node
4 transfers the exclusive lock on the data item to another node.

1 ~~32~~33. The method of Claim ~~31~~32 wherein the node changes the mode from local to global
2 prior to transferring the exclusive lock on the data item to another node.

1 ~~33~~34. The method of Claim ~~31~~32 wherein the node informs the master when the node has
2 completed writing the data item to persistent storage.

1 ~~34~~35. The method of Claim ~~34~~32 wherein the master informs the other node that the node
2 has completed writing the data item to persistent storage in response to the node
3 informing the master that the node has completed writing the data item to persistent
4 storage.

1 ~~35~~36. The method of Claim ~~32~~33 wherein the other node changes the mode from global to
2 local in response to a message from the master after the node has informed the master
3 that the node has completed writing the data item to persistent storage.

1 ~~36~~37. A method for managing a data item, the method comprising the steps of:
2 when a data item is transferred from one node to another node, performing the steps
3 of
4 if the data item has been dirtied by the node and a mode associated with the
5 data item is local, then changing the mode from local to global prior to
6 sending the data item to another node;
7 if the data item has not been dirtied by the node and the mode associated with
8 the data item is local, then sending the data item to the other node
9 without changing the mode;
10 allowing the other node to write the data item to persistent storage without
11 requesting permission if the mode is local; and
12 requiring the other node to obtain permission to write the data item to
13 persistent storage if the mode is global.

1 ~~37~~38. The method of Claim ~~36~~37 wherein the node transfers the data item to the other node
2 prior to completion of the node writing the data item to persistent storage.

1 3839. The method of Claim 3738 wherein, after completion of the node writing the data
2 item to persistent storage, the node sends a message to a master of the data item to
3 indicate that the data item has been written to persistent storage.

1 3940. The method of Claim 3839 wherein:
2 the other node receives the data item in global mode; and
3 the other node sends a request to the master of the node for permission to write the
4 data item; and
5 the master responds to said request by informing said other node to change said mode
6 from global to local.

1 4041. A method for managing versions of a data item, the method comprising the steps of:
2 when a dirty version of a data item is transferred from a first node to a second node
3 while a being-written version of the data item is being written to persistent
4 storage, performing the steps of:
5 communicating version information about the being-written version to the
6 second node; and
7 based on the version information, the second node preventing any version of
8 the data item that belongs to a first set of versions from being merged
9 with any version of the data item that belongs to a second set of
10 versions;
11 wherein the first set of versions includes all versions of the data item within
12 the second node that are at least as old as the being-written version;
13 and
14 wherein the second set of versions includes versions of the data item within
15 the second node that are newer than the being-written version.

1 4142. The method of Claim 4041 wherein the step of communicating is performed by a
2 master assigned to said data item.

1 4243. The method of Claim 4041 wherein:
2 the second node includes a plurality of versions in said first set; and
3 the second node merges said plurality of versions.

1 4344. The method of Claim 4041 further comprising the steps of:
2 informing the second node when the being-written version has been successfully
3 written to persistent storage; and
4 after the second node has been informed that the being-written version has been
5 successfully written to persistent storage, allowing said second node to discard
6 all versions in said first set of versions.

1 4445. The method of Claim 4243 further comprising the steps of:
2 informing the second node when the being-written version has been successfully
3 written to persistent storage; and
4 after the second node has been informed that the being-written version has been
5 successfully written to persistent storage, allowing said second node to discard
6 a merged version created by merging said plurality of versions.

1 4546. A method for managing past images of a data item, the method comprising the steps
2 of:
3 estimating a likelihood that a first past version of a data item will soon be written to
4 persistent storage or covered by a write to persistent storage;

5 if the estimated likelihood is exceeds a particular threshold, then storing a second past
6 version of the data item separate from the first past version of the data item;
7 and
8 if the estimated likelihood falls below a particular threshold, then merging the second
9 past version of the data item with the first past version of the data item.

1 4647. The method of Claim 4546 wherein the step of estimating is based on a comparison
2 between a time associated with the first past version of the data item and a time
3 associated with a recent entry in a redo log file.

1 4748. The method of Claim 4546 wherein the step of estimating is based on a comparison
2 between a time associated with the first past version of the data item and a time
3 associated with an entry at the head of a checkpoint queue.

1 4849. A computer-readable medium carrying instructions for managing caches in a system
2 with multiple caches that may contain different copies of a data item, the instructions
3 comprising instructions for performing the steps of:
4 modifying the data item in a first node of said multiple caches to create a modified
5 data item;
6 sending the modified data item from said first node to a second node of said multiple
7 caches without durably storing the modified data item from said first node to
8 persistent storage;
9 after said modified data item has been sent from said first node to said second node,
10 said first node sending a request to a master of said data item for writing said
11 data item to persistent storage; and
12 in response to said request, said master coordinating with said multiple caches to
13 cause said data item to be written to persistent storage.

1 4950. The computer-readable medium of Claim 4849 wherein:
2 the computer-readable medium includes instructions for performing the step of
3 maintaining, within an ordered series of bins, entries for past-image versions
4 of data items;
5 each bin in said ordered series corresponds to a time range;
6 a particular bin corresponds to the time range that covers the time at which the data
7 item is modified in said first node; and
8 the step of sending a request is performed by sending a request for writing a particular
9 bin of said ordered series of bins to persistent storage.

1 5051. The computer-readable medium of Claim 4950 wherein the step of said master
2 coordinating with said multiple caches to cause said data item to be written to
3 persistent storage includes said master causing said multiple caches to write data
4 items to persistent storage to cover all past image versions of data items that were
5 modified during the time range of said particular bin.

1 5152. The computer-readable medium of Claim 5051 further comprising instructions for
2 performing the step of emptying said particular bin after said multiple caches write
3 data items to persistent storage to cover all past image versions of data items that
4 were modified during the time range of said particular bin.

1 5253. The computer-readable medium of Claim 5152 wherein the step of emptying said
2 particular bin includes the steps of:
3 discarding entries within said particular bin that are associated with past images that
4 have last-dirtied times within the time range of said particular bin; and

5 moving to one or more other bins the entries within said particular bin that are
6 associated with past images that have last-dirtied times later than the time
7 range of said particular bin.

1 ~~53~~54. The computer-readable medium of Claim ~~48~~49 wherein the step of sending a request
2 to a master is performed by sending the request to a global lock manager.

1 ~~54~~55. The computer-readable medium of Claim ~~48~~49 wherein the step of sending a request
2 to a master is performed by sending the request to a lock manager that is one of a
3 plurality of lock managers within a distributed lock management system.

1 ~~55~~56. The computer-readable medium of Claim ~~48~~49 further comprising instructions for
2 performing the step of sending from the master, to interested nodes, write-notification
3 messages indicating that said data item has been written to persistent storage, in
4 response to said data item being written to persistent storage.

1 ~~56~~57. The computer-readable medium of Claim ~~55~~56 wherein the step of sending write-
2 notification messages includes the master sending to at least one interested node a
3 single message that notifies said at least one interested node that a plurality of data
4 items have been written to persistent storage.

1 ~~57~~58. The computer-readable medium of Claim ~~48~~49 wherein the step of said first node
2 sending a request to a master of said data item for writing said data item to persistent
3 storage includes the first node sending to said master a single message that requests
4 writing a plurality of data items to persistent storage, wherein said plurality of data
5 items includes said data item.

1 ~~58~~59. The computer-readable medium of Claim ~~57~~58 wherein the step of sending a single
2 message includes sending a message that identifies a bin to request that all data items
3 that belong to the bin be written to persistent storage.

1 ~~59~~60. The computer-readable medium of Claim ~~58~~59 wherein the bin is associated with a
2 range of time and includes data items that were first dirtied by the first node during
3 said range of time and that were subsequently transferred to other nodes without first
4 being written to persistent storage.

1 ~~60~~61. The computer-readable medium of Claim ~~55~~56 wherein the step of sending from the
2 master to interested nodes write-notification messages includes the steps of:
3 immediately sending write-notification messages to a first set of interested nodes,
4 where said first set of interested nodes includes the interested nodes that have
5 requested said data item to be written to persistent storage; and
6 delaying the sending of write-notification messages to a second set of nodes, where
7 said second set of nodes includes interested nodes that do not belong to said
8 first set of interested nodes.

1 ~~61~~62. The computer-readable medium of Claim ~~55~~56 wherein the step of sending from the
2 master to interested nodes includes delaying the sending of write-notification
3 messages to at least one interested node.

1 ~~62~~63. The computer-readable medium of Claim ~~61~~62 wherein a write-notification message
2 is sent to the at least one interested node in response to a lock request made by said at
3 least one interested node.

1 6364. The computer-readable medium of Claim ~~61~~62 wherein a write-notification message
2 is sent to the at least one interested node in response to the at least one interested node
3 requesting that said data item be written to persistent storage.

1 6465. The computer-readable medium of Claim ~~61~~62 wherein a write-notification message
2 is sent to the at least one interested node within a ping request that the master sends to
3 the at least one interested node for the at least one interested node to transfer another
4 data item to another node.

1 6566. The computer-readable medium of Claim ~~48~~49 wherein the step of coordinating
2 includes the steps of:
3 determining whether a version of said data item, that is at least as recent as said
4 modified version, has already been written to persistent storage; and
5 if a version of said data item that is at least as recent as said modified version has
6 already been written to persistent storage, then sending a write-notification
7 message from said master to notify said first node that a version of said data
8 item that is at least as recent as said modified version has already been written
9 to persistent storage.

1 6667. The computer-readable medium of Claim ~~65~~66 wherein the step of coordinating
2 includes, if a version of said data item that is at least as recent as said modified
3 version has not already been written to persistent storage, then sending a write-
4 perform message from said master to grant permission for said modified version to be
5 written to persistent storage.

1 6768. The computer-readable medium of Claim ~~48~~49 wherein the step of
2 coordinating includes the steps of:

3 selecting a particular node of said multiple caches that has a particular version of said
4 data item, wherein said particular version is at least as recent as the modified
5 data item in said first node; and
6 causing said particular version of said data item to be written from said particular
7 node to persistent storage.

1 ~~68~~69. The computer-readable medium of Claim ~~67~~68 wherein the step of selecting a
2 particular node includes selecting the node, of said multiple caches, that has a most
3 recently modified version of said data item.

1 ~~69~~70. The computer-readable medium of Claim ~~67~~68 further comprising instructions for
2 performing the step of the master informing the first node that said data item has been
3 written to persistent storage in response to the master receiving confirmation that said
4 particular version of said data item has been written to persistent storage.

1 ~~70~~71. The computer-readable medium of Claim ~~67~~68 further comprising instructions for
2 performing the step of the master informing a set of caches that said data item has
3 been written to persistent storage in response to the master receiving confirmation
4 that said particular version of said data item has been written to persistent storage,
5 wherein said set of caches includes caches, other than said particular node, that
6 contain modified versions of said data item that are not more recent than said
7 particular version.

1 ~~71~~72. A computer-readable medium carrying instructions for managing caches in a system
2 with multiple caches that may contain different copies of a data item, the instructions
3 comprising instructions for performing the steps of:
4 modifying the data item in a first cache to create a modified data item;

5 in response to writing the modified data item to persistent storage, performing the
6 steps of:
7 a node associated with the first cache determining whether any other cache in
8 said multiple caches had created a dirty version of said data item; and
9 if any other cache in said multiple caches had created a dirty version of said
10 data item, then the node associated with the first cache informing a
11 master of said data item that said modified data item has been written
12 to persistent storage; and
13 if no other cache in said multiple caches had created a dirty version of said
14 data item, then the step of writing the modified data item is performed
15 without informing said master that said modified data item has been
16 written to persistent storage.

1 ~~72~~73. The computer-readable medium of Claim ~~71~~72 wherein the step of determining
2 whether any other cache in said multiple caches had created a dirty version of said
3 data item includes inspecting a global dirty flag associated with said data item.

1 ~~73~~74. The computer-readable medium of Claim ~~72~~73 wherein:
2 prior to modifying the data item in said first cache, the data item had been modified in
3 a second cache;
4 the data item is not persistently stored between being modified in said second cache
5 and being modified in said first cache; and
6 a node associated with the second cache causes said global dirty flag to be set to
7 indicate that said data item is globally dirty.

1 ~~74~~75. The computer-readable medium of Claim ~~73~~74 wherein the node associated with the
2 second cache causes said global dirty flag to be set in response to transferring said
3 data item from said second cache to another cache of said multiple caches.

1 ~~75~~76. A computer-readable medium carrying instructions for managing caches in a system
2 with multiple caches that may contain different copies of a data item, the instructions
3 comprising instructions for performing the steps of:

4 modifying the data item in a first cache to create a modified data item;
5 when a node associated with the first node desires to write said modified data item to
6 persistent storage, performing the steps of:

7 if then node associated with the first cache does not currently have ownership
8 rights to said data item, then the node associated with the first cache
9 sending a request to a master of said data item for said data item to be
10 written to persistent storage; and

11 if said node associated with said first cache currently has ownership rights to
12 said data item, then the first node writing said data item to persistent
13 storage without sending a request to said master for said data item to
14 be written to persistent storage.

1 ~~75~~77. The computer-readable medium of Claim ~~74~~76 further comprising instructions for
2 performing the step of designating holders of exclusive locks in data items to be
3 owners of said data items.

1 ~~76~~78. A computer-readable medium carrying instructions for managing a data item, the
2 instructions comprising instructions for performing the steps of:
3 when a node that has an exclusive lock on a data item desires to write the data item to
4 persistent storage, performing the steps of

5 determining whether a mode associated with the data item is local or global;
6 if the mode associated with the data item is local, then the node writing the
7 data item to persistent storage without communicating with a master of
8 said data item; and
9 if the mode associated with the data item is global, then the node sending a
10 message to the master of the data item to request writing of said data
11 item to persistent storage.

1 ~~77~~79. The computer-readable medium of Claim ~~76~~78 wherein:
2 the mode associated with the data item is global; and
3 the computer-readable medium further includes instructions for:
4 the node receiving permission from the master to write the data item to
5 persistent storage; and
6 after writing the data item to persistent storage, changing the mode from
7 global to local.

1 ~~78~~80. The computer-readable medium of Claim ~~76~~78 wherein:
2 the mode associated with the data item is local; and
3 before the node has completed writing of the data item to persistent storage, the node
4 transfers the exclusive lock on the data item to another node.

1 ~~79~~81. The computer-readable medium of Claim ~~78~~80 wherein the node changes the mode
2 from local to global prior to transferring the exclusive lock on the data item to another
3 node.

1 ~~80~~82. The computer-readable medium of Claim ~~78~~80 wherein the node informs the master
2 when the node has completed writing the data item to persistent storage.

1 ~~81~~83. The computer-readable medium of Claim ~~78~~80 wherein the master informs the other
2 node that the node has completed writing the data item to persistent storage in
3 response to the node informing the master that the node has completed writing the
4 data item to persistent storage.

1 ~~82~~84. The computer-readable medium of Claim ~~79~~81 wherein the other node changes the
2 mode from global to local in response to a message from the master after the node has
3 informed the master that the node has completed writing the data item to persistent
4 storage.

1 ~~83~~85. A computer-readable medium carrying instructions for managing a data item, the
2 instructions comprising instructions for performing the steps of:
3 when a data item is transferred from one node to another node, performing the steps
4 of
5 if the data item has been dirtied by the node and a mode associated with the
6 data item is local, then changing the mode from local to global prior to
7 sending the data item to another node;
8 if the data item has not been dirtied by the node and the mode associated with
9 the data item is local, then sending the data item to the other node
10 without changing the mode;
11 allowing the other node to write the data item to persistent storage without
12 requesting permission if the mode is local; and
13 requiring the other node to obtain permission to write the data item to
14 persistent storage if the mode is global.

1 8486. The computer-readable medium of Claim 8385 wherein the node transfers the data
2 item to the other node prior to completion of the node writing the data item to
3 persistent storage.

1 ~~8587~~. The computer-readable medium of Claim 8486 wherein, after completion of the node
2 writing the data item to persistent storage, the node sends a message to a master of the
3 data item to indicate that the data item has been written to persistent storage.

1 ~~8688~~. The computer-readable medium of Claim 8587 wherein:
2 the other node receives the data item in global mode; and
3 the other node sends a request to the master of the node for permission to write the
4 data item; and
5 the master responds to said request by informing said other node to change said mode
6 from global to local.

1 8789. A computer-readable medium carrying instructions for managing versions of a data
2 item, the instructions comprising instructions for performing the steps of:
3 when a dirty version of a data item is transferred from a first node to a second node
4 while a being-written version of the data item is being written to persistent
5 storage, performing the steps of:
6 communicating version information about the being-written version to the
7 second node; and
8 based on the version information, the second node preventing any version of
9 the data item that belongs to a first set of versions from being merged
10 with any version of the data item that belongs to a second set of
11 versions;

12 wherein the first set of versions includes all versions of the data item within
13 the second node that are at least as old as the being-written version;
14 and
15 wherein the second set of versions includes versions of the data item within
16 the second node that are newer than the being-written version.

1 8890. The computer-readable medium of Claim 8789 wherein the step of communicating is
2 performed by a master assigned to said data item.

1 8991. The computer-readable medium of Claim 8789 wherein:
2 the second node includes a plurality of versions in said first set; and
3 the second node merges said plurality of versions.

1 9092. The computer-readable medium of Claim 8789 further comprising instructions for
2 performing the steps of:
3 informing the second node when the being-written version has been successfully
4 written to persistent storage; and
5 after the second node has been informed that the being-written version has been
6 successfully written to persistent storage, allowing said second node to discard
7 all versions in said first set of versions.

1 9193. The computer-readable medium of Claim 8991 further comprising instructions for
2 performing the steps of:
3 informing the second node when the being-written version has been successfully
4 written to persistent storage; and
5 after the second node has been informed that the being-written version has been
6 successfully written to persistent storage, allowing said second node to discard
7 a merged version created by merging said plurality of versions.

1 ~~92~~94. A computer-readable medium carrying instructions for managing past images of a
2 data item, the instructions comprising instructions for performing the steps of:
3 estimating a likelihood that a first past version of a data item will soon be written to
4 persistent storage or covered by a write to persistent storage;
5 if the estimated likelihood is exceeds a particular threshold, then storing a second past
6 version of the data item separate from the first past version of the data item;
7 and
8 if the estimated likelihood falls below a particular threshold, then merging the second
9 past version of the data item with the first past version of the data item.

1 ~~93~~95. The computer-readable medium of Claim ~~92~~94 wherein the step of estimating is
2 based on a comparison between a time associated with the first past version of the
3 data item and a time associated with a recent entry in a redo log file.

1 ~~94~~96. The computer-readable medium of Claim ~~92~~94 wherein the step of estimating is
2 based on a comparison between a time associated with the first past version of the
3 data item and a time associated with an entry at the head of a checkpoint queue.